# Multimodal Machine Learning

## Lecture 4.2: Aligned Representations

**Louis-Philippe Morency**

# Administrative Stuff

Language Technologies Institute

Carnegie Mellon University

# First Project Assignment

Due date: Sunday 9/24 at 8pm

Four main sections:

- Introduction
- Related work
- Experimental setup
- Research ideas

Follows ICML paper format

The two main sections are related work and research ideas

# teammates = # research ideas

Page limit depends on team size:
- 3 students : 4 pages + references
- 4 students : 4.5 pages + references
- 5 students : 5 pages + references

# Team Meetings with Instructor

- No lecture on Tuesday 10/3
- 15-mins meeting with instructor
  - Optional, but highly suggested
  - Not all teammates are required to attend
- Meetings next week: Wednesday 9/27 and Friday 9/29
- Signup form: https://calendly.com/morency/student-meetings

**Language Technologies Institute**

**Carnegie Mellon University**

# Multimodal Machine Learning

## Lecture 4.2: Aligned Representations

**Louis-Philippe Morency**

*\* Co-lecturer: Paul Liang. Original course co-developed with Tadas Baltrusaitis. Spring 2021 and 2022 editions taught by Yonatan Bisk.*
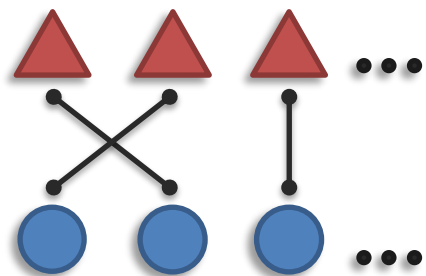
# Continuous Alignment

# Challenge 2: Alignment

**Definition:** Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure
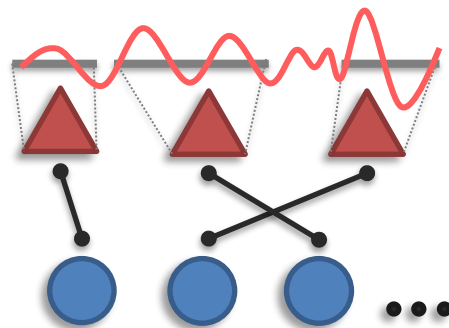
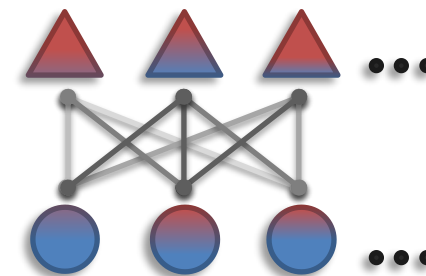## Sub-challenges:



**Discrete Alignment**

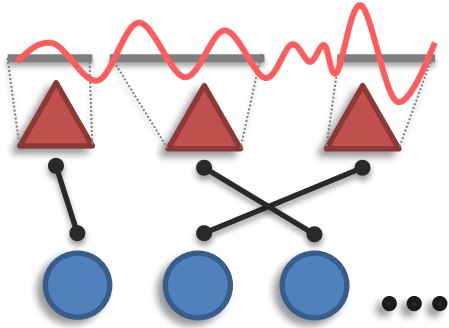Discrete elements and connections

**Continuous Alignment**

Segmentation and continuous warping
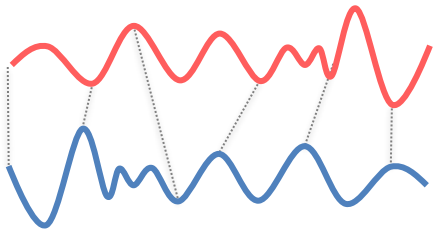
**Contextualized Representation**

Alignment + representation

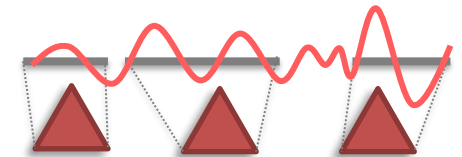# Challenge 2b: Continuous Alignment



**Definition:** Model alignment between modalities with continuous signals and no explicit elements
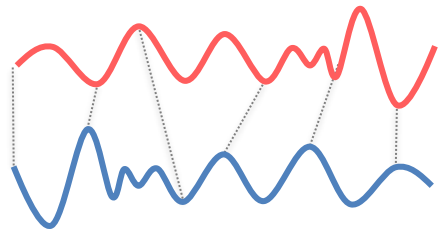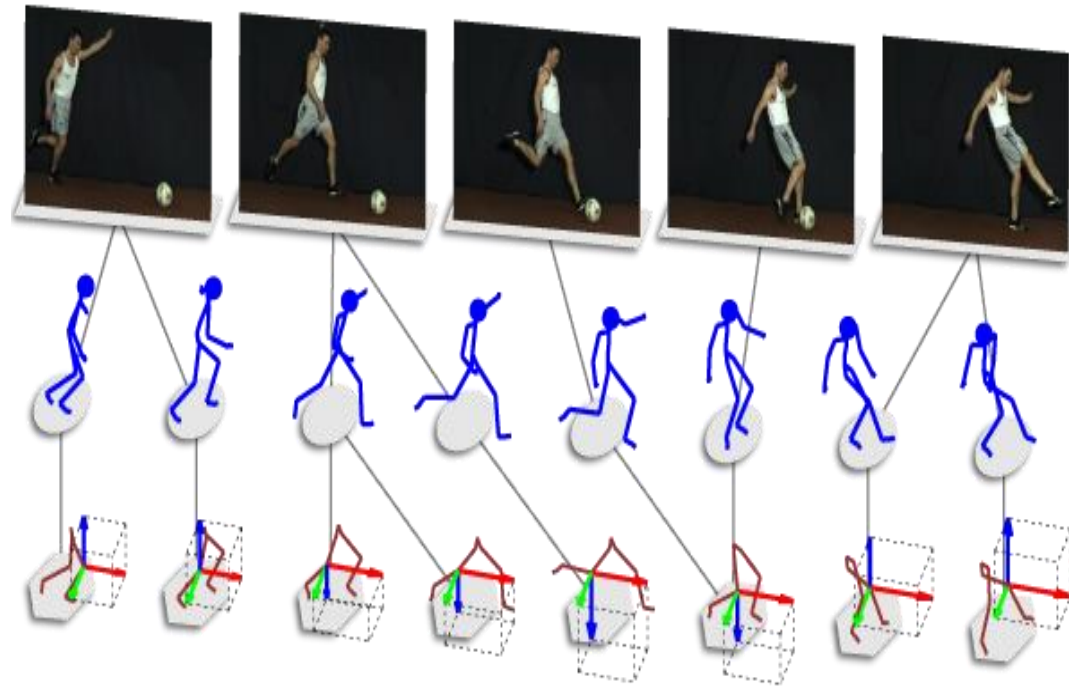
Continuous warping

Discretization
(segmentation)

# Continuous Warping – Example

Aligning video sequences

# Dynamic Time Warping (DTW)
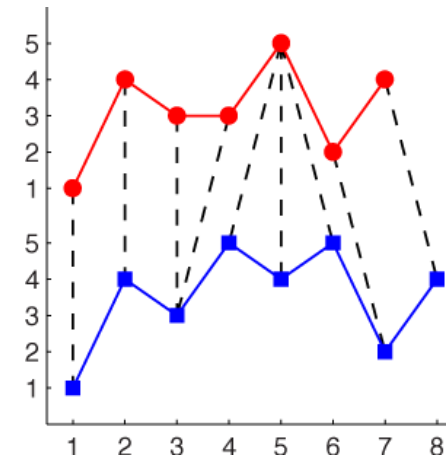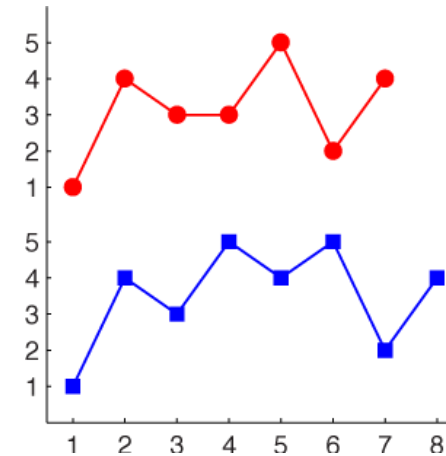
We have two unaligned temporal unimodal signals

- $\mathbf{X} = \left[ \boldsymbol{x}_1, \boldsymbol{x}_2, \ldots, \boldsymbol{x}_{n_x} \right] \in \mathbb{R}^{d \times n_x}$

- $\mathbf{Y} = \left[ \boldsymbol{y}_1, \boldsymbol{y}_2, \ldots, \boldsymbol{y}_{n_y} \right] \in \mathbb{R}^{d \times n_y}$

Find set of indices to minimize the alignment difference:

$$L(\boldsymbol{p}^x, \boldsymbol{p}^y) = \sum_{t=1}^{l} \left\| \boldsymbol{x}_{\boldsymbol{p}_t^x} - \boldsymbol{y}_{\boldsymbol{p}_t^y} \right\|_2^2$$

where $\boldsymbol{p}^x$ and $\boldsymbol{p}^y$ are index vectors of same length

Dynamic Time Warping is designed to find these index vectors!

# Dynamic Time Warping (DTW)

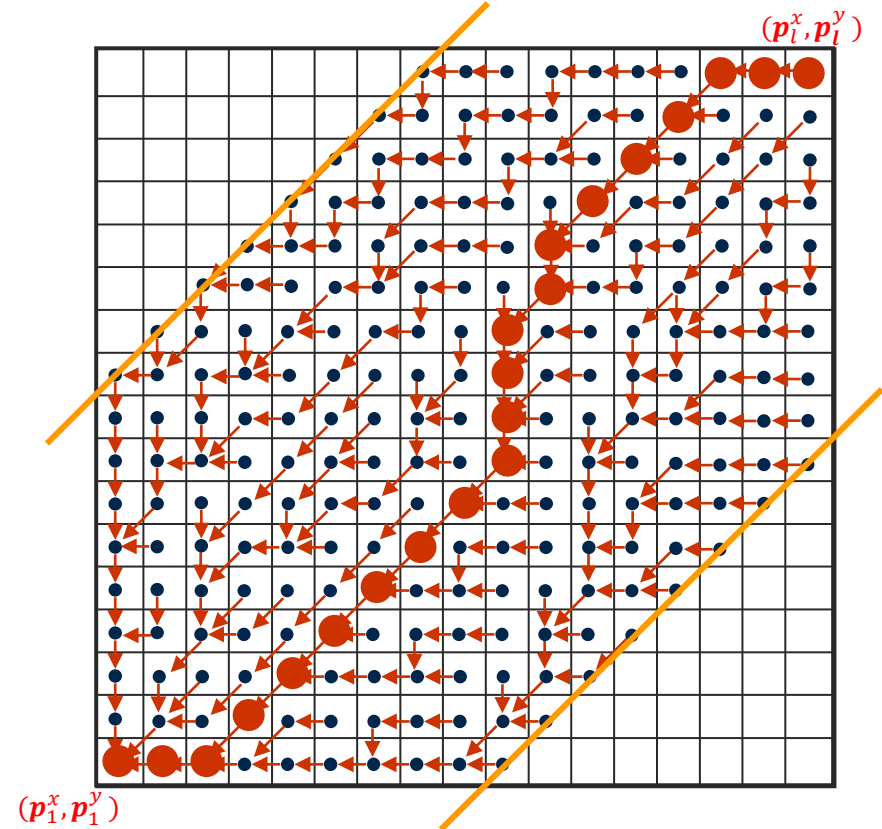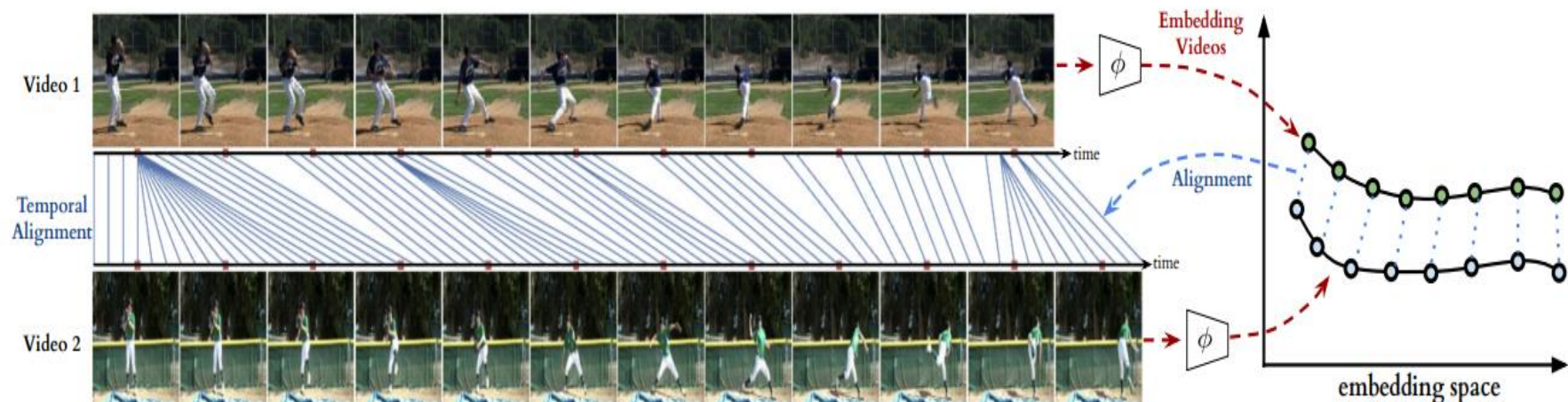Lowest cost path in a cost matrix

- Restrictions?
    - Monotonicity – no going back in time
    - Continuity  - no gaps
    - Boundary conditions - start and end at the same points
    - Warping window - don't get too far from diagonal
    - Slope constraint – do not insert or skip too much

Solved using dynamic programming while respecting the restrictions



$(p_l^x, p_l^y)$

$(p_1^x, p_1^y)$

# Temporal Alignment and Neural Representation Learning

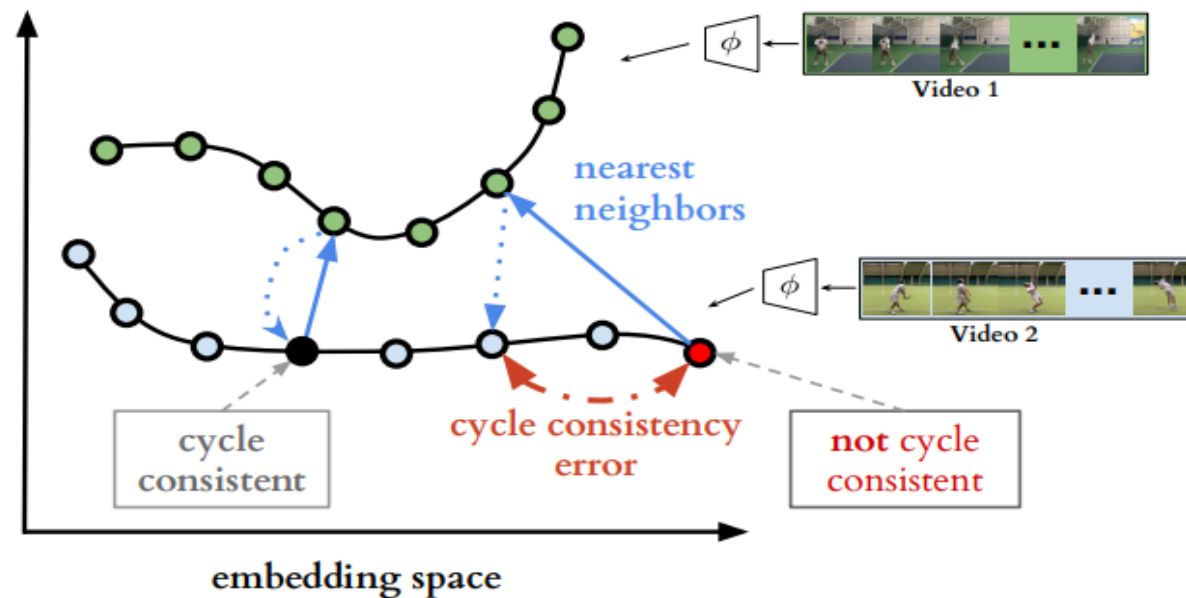**Premise:** we have paired video sequences
that can be be temporally aligned



How can we define a loss function to enforce
the alignment between sequences while at the
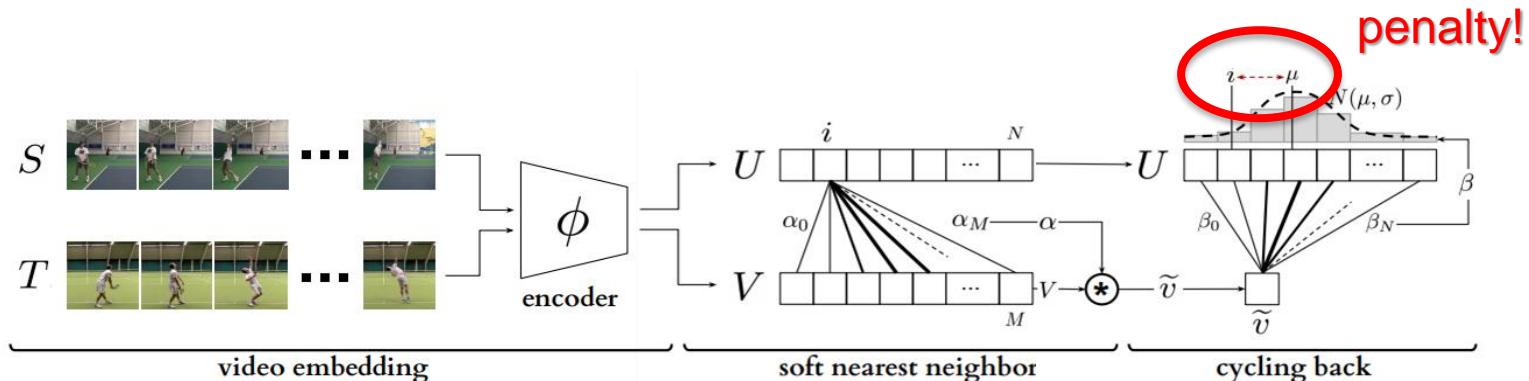same time learning good representations?

# Temporal Cycle-Consistency Learning

Solution: Representation learning by enforcing **Cycle consistency**



**Main idea:** My closest neighbor also views me as their closest neighbor

# Temporal Cycle-Consistency Learning



penalty!

video embedding  |  soft nearest neighbor  |  cycling back

Compute "soft" / "weighted" nearest neighbour:

distances: $\alpha_j = \dfrac{e^{-||u_i - v_j||^2}}{\sum_k^M e^{-||u_i - v_k||^2}}$

Soft nearest neighbor: $\widetilde{v} = \sum_j^M \alpha_j v_j,$

Find the nearest neighbor the other way and then penalize the distance:

$$\beta_k = \frac{e^{-||\widetilde{v} - u_k||^2}}{\sum_j^N e^{-||\widetilde{v} - u_j||^2}}$$

$$L_{cbr} = \frac{|i - \mu|^2}{\sigma^2} + \lambda \log(\sigma)$$

# Discretization (aka Segmentation)



**Common assumptions:** ① Segmented elements

Examples:



Medical imaging



Signals



Images

objects

# Discretization – Example

## Sequence Labeling and Alignment

**Phonemes**

| t | ah | m | aa | t | ow |
|---|----|---|----|----|----|

**Spectogram**

### How can we predict the sequence of phoneme labels?

Carnegie Mellon University

# Discretization – Example

## Sequence Labeling and Alignment

**Phonemes**

| t | ah | m | aa | t | ow |

**Spectogram**

**How can we predict the sequence of phoneme labels?**

**Challenge: many-to-1 alignment**

t    ah      m    aa

# Discretization – A Classification Approach

## Connectionist Temporal Classification

④    Most probable sequence labels

③    Predicted labels $l$

②    Path $\pi$ over the activations:

①    Output activations (distribution):

for 'blank' or no label   $y_{L+1}^t$

$y_L^t$

$y_1^t$

softmax

$l$

**Phonemes ($z$)**

t   ah   m   aa   t   ow

CTC

**Spectogram ($x$)**

Grave et al., Connectionist Temporal Classification: Labelling Unsegmented Sequence Data with Recurrent Neural Networks, ICML 2006

Language Technologies Institute     Carnegie Mellon University

**HUBERT: Hidden-Unit BERT**



K-mean clustering

# Challenge 2: Alignment

**Definition:** Identifying and modeling cross-modal connections between all elements of multiple modalities, building from the data structure

## Sub-challenges:



**Discrete Alignment**

Discrete elements and connections

**Continuous Alignment**

Segmentation and continuous warping

**Contextualized Representation**

Alignment + representation

# Contextualized Sequence Representations

# Sequence Encoding - Contextualization



**Option 1: Bi-directional LSTM:**
(e.g., ELMO)

How to encode this sequence while modeling the interaction between elements (e.g., words)?

But harder to parallelize…

# Sequence Encoding - Contextualization

**Option 2: Convolutions**



Can be parallelized!

But modeling long-range dependencies require multiple layers

And convolutional kernels are static

# Sequence Encoding - Contextualization

**Option 3: Self-attention**

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

Contextualized
Sequence Encoding

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I    do    not    like    it

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

**Self-Attention**

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I    do    not    like    it

Can be parallelized!

Long-range dependencies

Dynamic attention weights

# Self-Attention

# Self-Attention

# Self-Attention

Carnegie Mellon University

# Transformer Self-Attention

# Transformer Self-Attention

# Transformer Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer Self-Attention

Language Technologies Institute

Carnegie Mellon University

# Transformer Self-Attention

What if we want to attend simultaneously to multiple subspaces of $x$?

$h_1$ $h_2$ $h_3$ $h_4$ $h_5$

## Transformer's Self-Attention Layer

$W_q$ $W_k$ $W_v$

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$

I do not like it

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention



What happens if the words are shuffled?

# Position embeddings

❑ Position information is not encoded in a self-attention module

## How can we encode position information?

**Simple approach:** one-hot encoding

# Position embeddings

❑ Position information is not encoded in a self-attention module

How can we encode position information?

Simple approach: one-hot encoding + linear embeddings + $\begin{cases} \text{Sum} \\ \text{- or -} \\ \text{concat} \end{cases}$



| | | | | |
|---|---|---|---|---|
| $x_1$ $p_1$ | $x_2$ $p_2$ | $x_3$ $p_3$ | $x_4$ $p_4$ | $x_5$ $p_5$ |
| not | like | I | it | do |

# Transformer Multi-Head Self-Attention

# Transformer Multi-Head Self-Attention

In vector format…



$h$

Transformer's Multi-Head
Self-Attention Layer

$W_q^3$ $W_k^3$ $W_v^3$

$W_q^2$ $W_k^2$ $W_v^2$

$W_q^1$ $W_k^1$ $W_v^1$

$p$

$x$

# Transformer Multi-Head Attention

# Transformer – Residual Connection

# Language Pre-training

# Token-level and Sentence-level Embeddings

Token-level embeddings

Sentence-level embedding

# Pre-Training and Fine-Tuning



**Pre-training**
(e.g., language model)

**Fine-Tuning**

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

①    Jointly learn representation for token-level and sentence level

②    Same network architecture for pre-training and fine-tuning

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

③ Can be used learn relationship between sentences

④ Models bidirectional and long-range interactions between tokens

How can we do all this?

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

| $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I      do     not    like   it            I      enjoy   my     time   here

# BERT: Bidirectional Encoder Representations from Transformers

**Advantages:**

① Jointly learn representation for token-level and sentence level

② Same network architecture for pre-training and fine-tuning

③ Can be used learn relationship between sentences

④ Models bidirectional interactions between tokens

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

Special sentence-level token

But how to train self-supervised?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I   do   not   like   it       I   enjoy   my   time   here

# Pre-training BERT Model

**1** **Masked Language Model**

Randomly mask input tokens and then try to predict them

What is the loss function?

# Pre-training BERT Model

② **Next Sentence Prediction**

Given two sentences, predict if this is the next one or not

What is the loss function?

Where can we find training data?

How can BERT know the difference between both sentences?

IsNext
- or -
NotNext

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

I    do    not    like    it        I    enjoy    my    time    here

# Three Embeddings: Token + Position + Sentence



| Input | [CLS] | my | dog | is | cute | [SEP] | he | likes | play | ##ing | [SEP] |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Token Embeddings | $E_{[CLS]}$ | $E_{my}$ | $E_{dog}$ | $E_{is}$ | $E_{cute}$ | $E_{[SEP]}$ | $E_{he}$ | $E_{likes}$ | $E_{play}$ | $E_{\#\#ing}$ | $E_{[SEP]}$ |
| Segment Embeddings | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_A$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ | $E_B$ |
| Position Embeddings | $E_0$ | $E_1$ | $E_2$ | $E_3$ | $E_4$ | $E_5$ | $E_6$ | $E_7$ | $E_8$ | $E_9$ | $E_{10}$ |

# Fine-Tuning BERT

① Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**How?**



| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

|  | I | do | not | like | it |  | I | enjoy | my | time | here |

# Fine-Tuning BERT

Sentence-level classification for only one sentence

Examples: sentiment analysis, document classification

**And if we have a label for each token?**

# Fine-Tuning BERT

② Token-level classification for only one sentence

Examples: part-of-speech tagging, slot filling



How to compare two sentences?

# Fine-Tuning BERT

Examples: natural language inference

$$\hat{y}_s$$

| softmax |
|---------|

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

|   | I | do | not | like | it | | I | enjoy | my | time | here |

# Fine-Tuning BERT

④ Question-answering: find start/end of the answer in the document

**Paragraph:** " ... *Other legislation followed, including the Migratory Bird Conservation Act of 1929, a* 1937 treaty *prohibiting the hunting of right and gray whales, and the* Bald Eagle Protection Act of 1940. *These* later laws *had a low cost to society—the species were relatively rare—and little* opposition *was raised.*"

**Question 1:** "*Which laws faced significant* opposition?*"
**Plausible Answer:** later laws

**Question 2:** "*What was the name of the* 1937 treaty?*"
**Plausible Answer:** Bald Eagle Protection Act

| $h_s$ | $h_1$ | $h_2$ | $h_3$ | $h_4$ | $h_5$ | $h_{sep}$ | $h'_1$ | $h'_2$ | $h'_3$ | $h'_4$ | $h'_5$ |

## Transformer Self-Attention

How?

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

How  old  is  the  man      He  is  25  years  old

# Fine-Tuning BERT

④ Question-answering: find start/end of the answer in the document



Maximum value gives start time

Same architecture for the end time

softmax

Learned during fine-tuning

**Transformer Self-Attention**

| cls | $x_1$ | $x_2$ | $x_3$ | $x_4$ | $x_5$ | sep | $x'_1$ | $x'_2$ | $x'_3$ | $x'_4$ | $x'_5$ |

How    old    is    the    man        He    is    25    years    old

# Sequence-to-Sequence Using Transformer

# Sequence-to-Sequence Modeling

Je      n'      aime      pas      cela

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$

## How can we perform seq2seq translation with transformer attention?

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I      do      not      like      it

# Seq2Seq with Transformer Attentions

Je    n'    aime    pas    cela

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$

$h_1$    $h_2$    $h_3$    $h_4$    $h_5$

self-attention

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I    do    not    like    it

# Seq2Seq with Transformer Attentions

Je    n'    aime    pas    cela

$\hat{y}_1$  $\hat{y}_2$  $\hat{y}_3$  $\hat{y}_4$  $\hat{y}_5$

$h_1$  $h_2$  $h_3$  $h_4$  $h_5$

self-attention

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$

I    do    not    like    it

$g_1$  $g_2$  $g_3$  $g_4$  $g_5$

"masked" self-attention

$y_0$  $y_1$  $y_2$  $y_3$  $y_4$

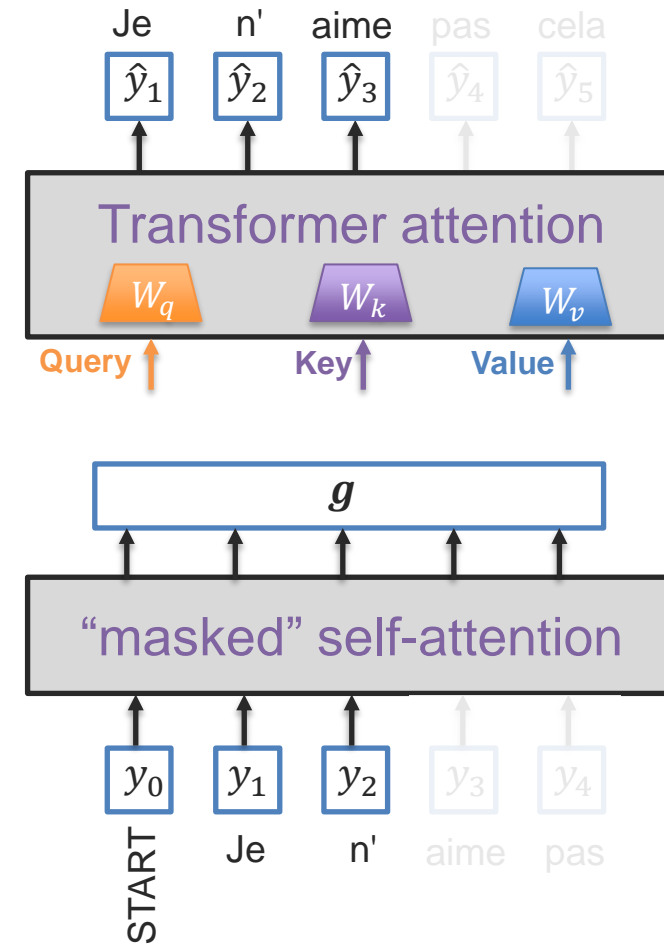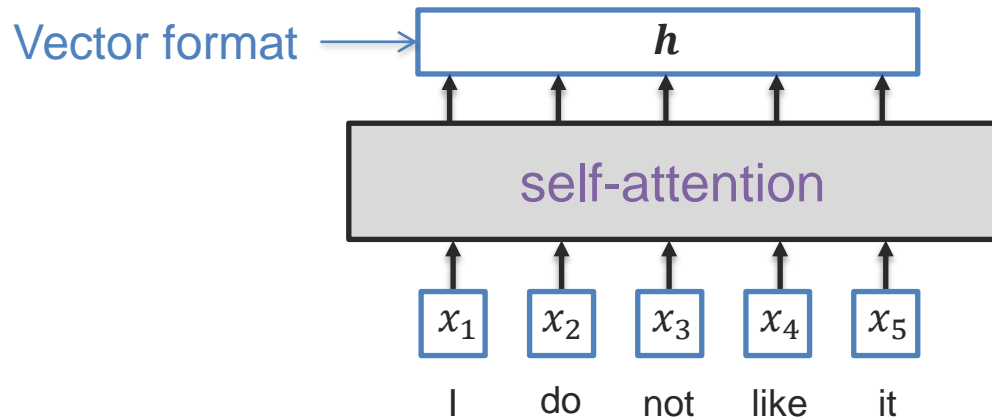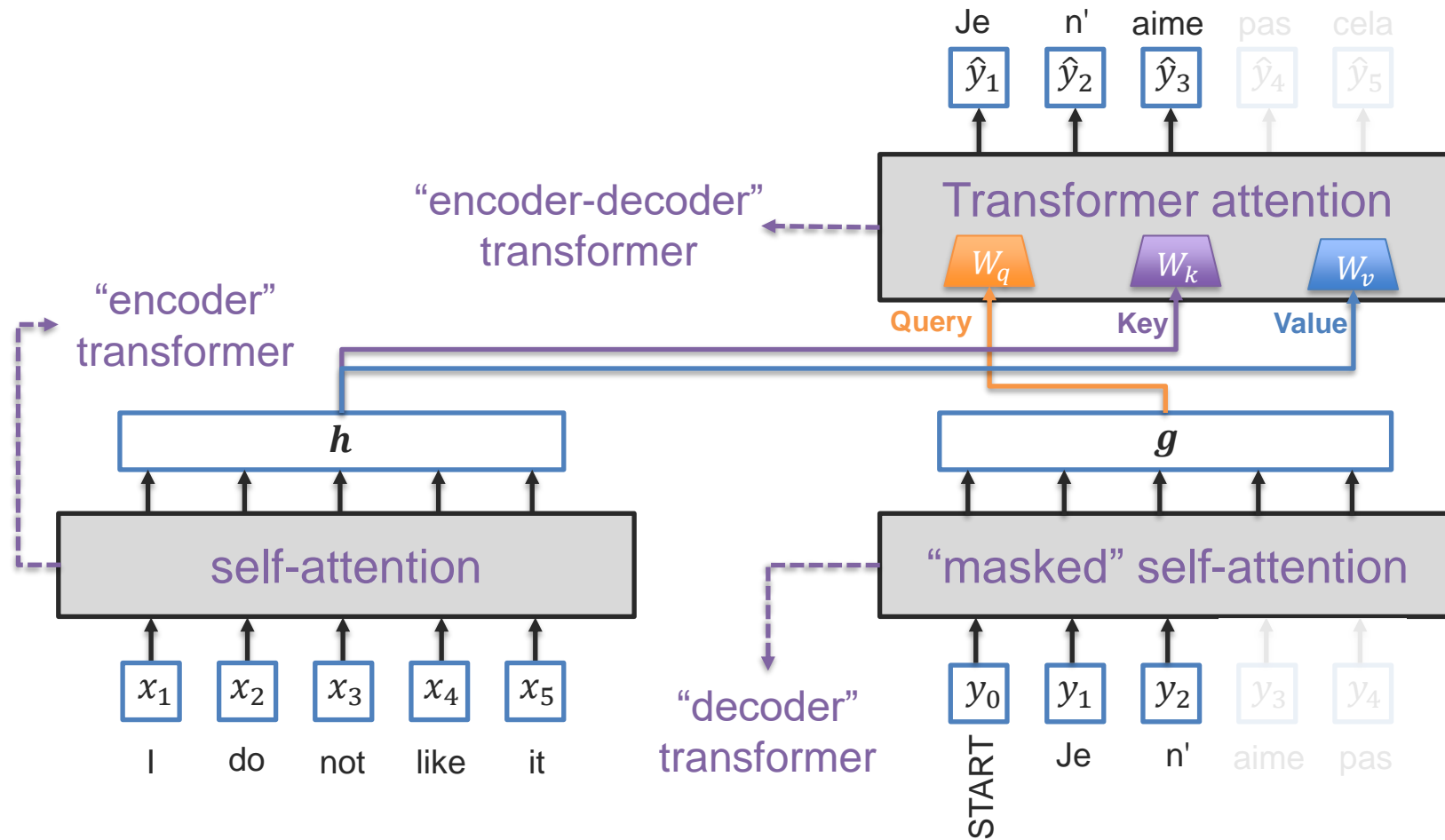START    Je    n'    aime    pas

# Seq2Seq with Transformer Attentions

How should we connect the encoder and decoder self-attention to the transformer attention?

Je    n'    aime    pas    cela

$\hat{y}_1$    $\hat{y}_2$    $\hat{y}_3$    $\hat{y}_4$    $\hat{y}_5$

Transformer attention

$W_q$    $W_k$    $W_v$

Query    Key    Value

Vector format ⟶    $h$

$g$

self-attention

"masked" self-attention

$x_1$    $x_2$    $x_3$    $x_4$    $x_5$

I    do    not    like    it

$y_0$    $y_1$    $y_2$    $y_3$    $y_4$

START    Je    n'    aime    pas

# Seq2Seq with Transformer Attentions

# And Many More… Next week!